

Bayesian Adaptive Lasso with Variational Bayes for Variable Selection in High-dimensional Generalized Linear Mixed Models

Dao Thanh Tung* Minh-Ngoc Tran[†] Tran Manh Cuong*

August 31, 2016

Abstract

This article describes a full Bayesian treatment for simultaneous fixed-effect selection and parameter estimation in high-dimensional generalized linear mixed models. The approach consists of using a Bayesian adaptive Lasso penalty for signal-level adaptive shrinkage and a fast Variational Bayes scheme for estimating the posterior mode of the coefficients. The proposed approach offers several advantages over the existing methods, for example, the adaptive shrinkage parameters are automatically incorporated, no Laplace approximation step is required to integrate out the random effects. The performance of our approach is illustrated on several simulated and real data examples. The algorithm is implemented in the R package `glmmbv` and is made available online.

Keywords: Posterior mode, Lasso, High dimensions, EM algorithm

1 Introduction

Generalized linear mixed models (GLMMs) are widely used for modeling cluster-dependent data. Variable selection in GLMMs is considered a difficult task, because of the present of integrals that are often analytically intractable. Classical methods for variable selection, such as the ones based on hypothesis testing or subset selection, are restricted to a few covariates. Notable works are two recent papers by Groll and Tutz (2012) and Schelldorfer et al. (2013) which can do variable selection for GLMMs in high dimensions. Their approach first estimates the likelihood by approximating the integrals over the random effects using the Laplace method, then minimizes the sum of this estimated likelihood and a Lasso-type penalty which is the l_1 -norm of the fixed effect coefficients. Using a Lasso-type penalty will shrink the coefficients towards zero, thus leading to variable selection. This variable selection approach is attractive compared to the classical approaches as it can handle problems with a large number of potential covariates.

However, there is still room for improvement within the approach of Groll and Tutz (2012) and Schelldorfer et al. (2013). First, the Laplace approximation of the likelihood might be in

*Vietnam National University, Hanoi

[†]The University of Sydney Business School. Correspondence to: minh-ngoc.tran@sydney.edu.au

some cases not very accurate (see, e.g. Joe, 2008). Second, the performance depends on the shrinkage parameter that needs to be selected appropriately. So that the user has to run the procedure over and over again for different values of the shrinkage parameter within a pre-specified range, then selects the best value of the shrinkage parameter based on some criterion such as AIC or BIC. As the result, the entire procedure for selecting the final model may be time consuming. Furthermore, specifying an appropriate range for the shrinkage parameter is not straightforward. Third, this approach uses the same shrinkage parameter for every coefficients, which can lead to biased estimates of the coefficients.

This article proposes using the Bayesian adaptive Lasso for variable selection in high-dimensional GLMMs. We use double exponential priors for the coefficients with different shrinkage parameters for different coefficients, which is equivalent to the approach in Groll and Tutz (2012) and Schelldorfer et al. (2013) when all the shrinkage parameters are equal. It is desirable to apply different shrinkage on different coefficients to achieve adaptivity, i.e. larger shrinkage should be put on coefficients corresponding to unimportant covariates and smaller shrinkage should be used for important covariates (Zou, 2006). We consider a full Bayesian treatment, i.e. we put appropriate priors on all the model parameters, including the shrinkage parameters. As the result, we overcome the challenging task of selecting a high-dimensional vector of the shrinkage parameters.

We then develop a variational Bayes (VB) algorithm for estimating the posterior mode of the coefficient vector and the posterior distribution of the covariance matrix of the random effects. This leads to a totally automatic procedure for simultaneous variable selection and parameter estimation in GLMMs, and the adaptive shrinkage parameters are automatically incorporated. Finally, unlike the approach in Groll and Tutz (2012) and Schelldorfer et al. (2013), our approach does not rely on the Laplace approximation for integrating out the random effects, because the updating procedure in the variational Bayes algorithm leads to an integral that either can be computed analytically or approximated in close form with an arbitrary accuracy. The examples in Section 4 show that our approach outperforms the existing methods in terms of the rate of correctly-fitted models, the mean squared error of the estimates, and the CPU running time.

The paper is organized as follows. Section 2 provides some background on the variational Bayes method, and presents the VB method for estimating the posterior mode. Section 3 describes our algorithm for variable selection in GLMMs. Section 4 presents a systematic simulation example and real data applications. Section 5 concludes and discusses some possible extensions. The algorithm is implemented in the R package `glmmb` and is available at <https://sites.google.com/site/mntran26/research>.

2 Variational Bayes method

Suppose we have data y , a likelihood $p(y|\theta)$ where $\theta \in \mathbb{R}^d$ is an unknown parameter, and a prior distribution $p(\theta)$ for θ . Variational Bayes (VB) approximates the posterior $p(\theta|y) \propto p(\theta)p(y|\theta)$ by a distribution $q(\theta)$ within some more tractable class, chosen to minimize the Kullback-Leibler divergence

$$\text{KL}(q||p) = \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta. \quad (1)$$

We have

$$\log p(y) = \int q(\theta) \log \frac{p(y, \theta)}{q(\theta)} d\theta + \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta = L(q) + \text{KL}(q\|p),$$

where

$$L(q) = \int q(\theta) \log \frac{p(y, \theta)}{q(\theta)} d\theta. \quad (2)$$

As $\text{KL}(q\|p) \geq 0$, $\log p(y) \geq L(q)$ for every $q(\theta)$. $L(q)$ is therefore often called the lower bound, and minimizing $\text{KL}(q\|p)$ is equivalent to maximizing $L(q)$.

Often factorized approximations to the posterior are considered in variational Bayes. We explain the idea for a factorization with 2 blocks. Assume that $\theta = (\theta_1, \theta_2)$ and that $q(\theta)$ is factorized as

$$q(\theta) = q_1(\theta_1)q_2(\theta_2). \quad (3)$$

We further assume that $q_1(\theta_1) = q_{\tau_1}(\theta_1)$ and $q_2(\theta_2) = q_{\tau_2}(\theta_2)$ where τ_1 and τ_2 are variational parameters that need to be estimated. Then

$$\begin{aligned} L(\tau_1, \tau_2) = L(q) &= \int q_{\tau_1}(\theta_1)q_{\tau_2}(\theta_2) \log p(y, \theta) d\theta_1 d\theta_2 - \int q_{\tau_1}(\theta_1) \log q_{\tau_1}(\theta_1) d\theta_1 + C(\tau_2) \\ &= \int q_{\tau_1}(\theta_1) \left(\int q_{\tau_2}(\theta_2) \log p(y, \theta) d\theta_2 \right) d\theta_1 - \int q_{\tau_1}(\theta_1) \log q_{\tau_1}(\theta_1) d\theta_1 + C(\tau_2) \\ &= \int q_{\tau_1}(\theta_1) \log \tilde{p}(y, \theta_1) d\theta_1 - \int q_{\tau_1}(\theta_1) \log q_{\tau_1}(\theta_1) d\theta_1 + C(\tau_2) \\ &= \int q_{\tau_1}(\theta_1) \log \frac{\tilde{p}_1(y, \theta_1)}{q_{\tau_1}(\theta_1)} d\theta_1 + C(\tau_2), \end{aligned}$$

where $C(\tau_2)$ is a constant depending only on τ_2 and

$$\tilde{p}_1(y, \theta_1) = \exp \left(\int q_{\tau_2}(\theta_2) \log p(y, \theta) d\theta_2 \right) = \exp (E_{-\theta_1}(\log p(y, \theta))).$$

Given that τ_2 is fixed. Let

$$\tau_1^* = \tau_1^*(\tau_2) = \arg \max_{\tau_1} \left\{ \int q_{\tau_1}(\theta_1) \log \frac{\tilde{p}_1(y, \theta_1)}{q_{\tau_1}(\theta_1)} d\theta_1 \right\}, \quad (4)$$

then

$$L(\tau_1^*, \tau_2) \geq L(\tau_1, \tau_2) \text{ for all } \tau_1. \quad (5)$$

Similarly, given a fixed τ_1 , let

$$\tau_2^* = \tau_2^*(\tau_1) = \arg \max_{\tau_2} \left\{ \int q_{\tau_2}(\theta_2) \log \frac{\tilde{p}_2(y, \theta_2)}{q_{\tau_2}(\theta_2)} d\theta_2 \right\}, \quad (6)$$

with

$$\tilde{p}_2(y, \theta_2) = \exp \left(\int q_{\tau_1}(\theta_1) \log p(y, \theta) d\theta_1 \right) = \exp (E_{-\theta_2}(\log p(y, \theta))).$$

Then,

$$L(\tau_1, \tau_2^*) \geq L(\tau_1, \tau_2) \text{ for all } \tau_2. \quad (7)$$

Let $\tau^{\text{old}} = (\tau_1^{\text{old}}, \tau_2^{\text{old}})$ be the current value of τ_1 and τ_2 , update $\tau_1^{\text{new}} = \tau_1^*(\tau_2^{\text{old}})$ as in (4) and $\tau_2^{\text{new}} = \tau_2^*(\tau_1^{\text{new}})$ as in (6). Then, because of (5) and (7),

$$L(\tau^{\text{new}}) \geq L(\tau^{\text{old}}). \quad (8)$$

This leads to an iterative scheme for updating τ and (8) ensures the improvement of the lower bound over the iterations. Because the lower bound $L(\tau)$ is bounded from above by $\log p(y)$, the convergence of the iterative scheme is guaranteed. The above argument can be easily extended to the general case in which $q(\theta)$ is factorized into K blocks $q(\theta) = q_1(\theta_1) \times \dots \times q_K(\theta_K)$.

The variational Bayes approximation is now reduced to solving an optimization problem in the form of (4). Let $\tilde{p}_1(\theta_1|y)$ be the density of θ_1 determined by the unnormalized function $\tilde{p}_1(y, \theta_1)$, i.e.

$$\tilde{p}_1(\theta_1|y) = \frac{\tilde{p}_1(y, \theta_1)}{\int \tilde{p}_1(y, \theta_1) d\theta_1} \propto \exp(E_{-\theta_1}(\log p(y, \theta))). \quad (9)$$

In many cases, a conjugate prior $p(\theta_1)$ can be selected such that $\tilde{p}_1(\theta_1|y)$ belongs to a family of recognizable parametric densities. Then the optimal VB posterior $q_{\tau_1^*}(\theta_1)$ that maximizes the integral on the right hand side of (4) is $\tilde{p}_1(\theta_1|y)$, with τ_1^* the corresponding parameter of this density.

If $\tilde{p}_1(\theta_1|y)$ does not belong to a recognizable density family, some optimization technique is needed to solve (4). Note that (4) has exactly the same form as the original VB problem that attempts to maximize $L(q)$ in (2). We can first select a functional form for the variational distribution q and then estimate the unknown parameters accordingly. If the variational distribution is assumed to belong to the exponential family with unknown parameters τ , Salimans and Knowles (2013) propose a stochastic approximation method for solving for τ . The reader is referred to their paper for the details.

2.1 Variational Bayes method for estimating the posterior mode

As pointed out in Tibshirani (1996), the Lasso estimator is equivalent to the posterior mode when a double-exponential prior (also called Laplace prior) is used for the vector of coefficients β . In general, for the variable selection purposes in Bayesian settings, one is interested in the posterior mode rather than the entire posterior distribution. As will be seen in the next section, variable selection in GLMMs is carried out through computing the posterior mode of the fixed-effect coefficient vector β . We will present in this section a Variational Bayes method for estimating a posterior mode.

Write $\theta = (\theta_1, \theta_2)$, where θ_1 is the vector of parameters whose posterior mode is of our interest, and θ_2 is a vector of other parameters, random effects or missing data. Then, we can use a VB posterior of the form

$$q(\theta) = \delta_{\tau_1}(\theta_1) q_{\tau_2}(\theta_2), \quad (10)$$

with $\delta_{\tau_1}(\theta_1)$ a point mass density concentrated at τ_1 . For our purposes, τ_1 will be the estimate of the posterior mode of θ_1 .

Equations (4) and (6) become

$$\tau_1^*(\tau_2) = \arg \max_{\tau_1} \int q_{\tau_2}(\theta_2) \log p(y, \tau_1, \theta_2) d\theta_2, \quad (4')$$

and

$$\tau_2^*(\tau_1) = \arg \max_{\tau_2} \left\{ \int q_{\tau_2}(\theta_2) \log \frac{p(y, \tau_1, \theta_2)}{q_{\tau_2}(\theta_2)} d\theta_2 \right\}. \quad (6')$$

The optimal VB posterior of θ_2 from (6') is $q_{\tau_2^*}(\theta_2) = p(\theta_2|y, \tau_1) \propto p(y, \tau_1, \theta_2)$. Then (4') and (6') can be written in terms of the EM algorithm (Dempster et al., 1977), where

- E-step: compute $Q(\tau_1|\tau_1^{\text{old}}) = \int p(\theta_2|y, \tau_1^{\text{old}}) \log p(y, \tau_1, \theta_2) d\theta_2$.
- M-step: maximize $Q(\tau_1|\tau_1^{\text{old}})$ over τ_1 .

The EM algorithm therefore can be considered as a special case of this VB algorithm where $q_{\tau_2}(\theta_2)$ in (10) is $q_{\tau_2}(\theta_2) = p(\theta_2|y, \tau_1)$. Note that the VB mode method in (4') and (6') is somewhat more flexible than the EM algorithm because we have more freedom to find a solution to (6') provided that $q_{\tau_2}(\theta_2)$ is restricted to some density family. This is important because the optimal density $q_{\tau_2^*}(\theta_2) = p(\theta_2|y, \tau_1)$ in some cases does not belong to a family of recognizable densities, and it is difficult to compute the integral in the E-step. For example, in generalized linear mixed models considered in this paper, the distribution of the random effects conditional on the data and the other parameters does not belong to a family of recognizable densities, making it difficult to estimate the coefficient vector using the EM algorithm.

3 Variable selection and estimation for GLMMs

Consider a generalized linear mixed model where $y_i = (y_{i1}, \dots, y_{in_i})'$ is the vector of responses for the i th subject, $i = 1, \dots, m$. Given random effects b_i , the y_{ij} are conditionally independently distributed with the density or probability function

$$f(y_{ij}|\beta, b_i) = \exp \left(\frac{y_{ij}\eta_{ij} - \zeta(\eta_{ij})}{\phi} + c(y_{ij}, \phi) \right),$$

where η_{ij} is a canonical parameter which is monotonically related to the conditional mean $\mu_{ij} = E(y_{ij}|\beta, b_i)$ through a link function $g(\cdot)$, $g(\mu_{ij}) = \eta_{ij}$. The fixed effect coefficient vector is $\beta = (\beta_0, \beta'_{1:p})'$ with β_0 the slope and $\beta_{1:p} = (\beta_1, \dots, \beta_p)'$. The scale parameter ϕ can be unknown and $\zeta(\cdot)$ and $c(\cdot)$ are known functions. Here, for simplicity, we are considering the case of a canonical link function, i.e. $g(\mu_{ij}) = \eta_{ij}$. The vector $\eta_i = (\eta_{i1}, \dots, \eta_{in_i})'$ is modeled as $\eta_i = \beta_0 \mathbf{1}_{n_i} + X_i \beta_{1:p} + Z_i b_i$, where $\mathbf{1}_{n_i}$ is the vector of ones, X_i is an $n_i \times p$ design matrix for the fixed effects and Z_i is an $n_i \times u$ design matrix for the random effects (where u is the dimension of b_i). Let $n = \sum_{i=1}^m n_i$, $b = (b'_1, \dots, b'_m)'$ and

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}, \quad X = \begin{pmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_m \end{pmatrix}, \quad Z = \begin{pmatrix} Z_1 & 0 & \cdots & 0 \\ 0 & Z_2 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & Z_m \end{pmatrix}, \quad \eta = \begin{pmatrix} \eta_1 \\ \eta_2 \\ \vdots \\ \eta_m \end{pmatrix} = X\beta + Zb.$$

The likelihood conditional on the random effects b is

$$p(y|\beta, b, \phi) = \prod_{i=1}^m \prod_{j=1}^{n_i} f(y_{ij}|\beta, b_i) = \exp \left(\frac{1}{\phi} (y'\eta - \mathbf{1}'\zeta(\eta)) + c(y, \phi) \right),$$

where $\zeta(\eta)$ is understood component-wise and $c(y, \phi) = \sum_{i,j} c(y_{ij}, \phi)$.

The random effects b_i are often assumed independently distributed as $\mathcal{N}(0, Q^{-1})$, where $\mathcal{N}(\mu, \Sigma)$ denotes the multivariate normal distribution with mean μ and covariance matrix Σ . The distribution of b is $\mathcal{N}(0, Q_b^{-1})$ with Q_b a block diagonal matrix $\text{diag}(Q, \dots, Q)$. We consider Bayesian inference with the following hierarchy

$$\begin{aligned} y|\beta, b, \phi &\sim p(y|\beta, b, \phi) \\ b|Q &\sim \mathcal{N}(0, Q_b^{-1}) \\ Q &\sim \text{Wishart}(S_0, \nu_0) \\ p(\beta_0) &\sim 1 \\ \beta_j|\lambda_j &\sim \text{DE}(\lambda_j) = \frac{\lambda_j}{2} \exp(-\lambda_j|\beta_j|), \quad j = 1, \dots, p \\ \lambda_j &\sim \text{Gamma}(r, s) = \frac{s^r}{\Gamma(r)} (\lambda_j)^{r-1} \exp(-s\lambda_j), \end{aligned} \tag{11}$$

where $\text{DE}(\lambda_j)$ denotes the double-exponential density. If ϕ is unknown we also put a prior $p(\phi)$ on ϕ . We refer to the suggested model (11) as the Bayesian adaptive Lasso model (BaLasso) for GLMM. The set of model parameters is $\theta = (\beta, Q, \phi, b, \lambda_1, \dots, \lambda_p)$ and S_0, ν_0, r, s are hyperparameters whose selection is discussed later.

When $\lambda_j = \lambda$ and considered fixed, the joint posterior distribution of β, Q, ϕ is

$$p(\beta, Q, \phi) \propto p(\phi)p(Q) \exp \left(\log \int p(y|\beta, b, \phi)p(b|Q)db - \lambda \sum_{j=1}^p |\beta_j| \right).$$

In this case, the posterior marginal mode of β from model (11) is exactly the penalized maximum likelihood estimate in Groll and Tutz (2012) and Schellhdorfer et al. (2013), who estimate the parameters by maximizing

$$\log \int p(y|\beta, b, \phi)p(b|Q)db - \lambda \sum_{j=1}^p |\beta_j| \tag{12}$$

over β . Note that we use different λ_j for different coefficient β_j to achieve signal-level adaptivity (Zou, 2006).

The Bayesian Lasso was first proposed in Park and Casella (2008) who considered a single shrinkage λ for all coefficients, in the context of ordinary linear regression only. The Bayesian adaptive Lasso for GLMs was proposed in Griffin and Brown (2011) and Leng et al. (2013). Griffin and Brown (2011) employed the EM algorithm to estimate the posterior mode of β and were therefore able to carry out variable selection. Leng et al. (2013) first used Gibbs sampling to sample from the posterior of λ and then proposed a Bayesian-frequentist hybrid method for doing variable selection where λ is fixed to its posterior mode. To the best of our knowledge, this paper is the first to consider the Bayesian adaptive Lasso model (11) for inference in GLMMs, and also use VB for estimating a posterior mode.

We use Variational Bayes to approximate the posterior $p(\theta|y)$ with the variational posterior factorized as

$$q(\theta) = q(\beta)q(Q)q(\phi)q(b) \prod_{j=1}^p q(\lambda_j) \tag{13}$$

where $q(\beta) = \delta_{\beta^q}(\beta)$ and $q(b)$ is normal with mean μ_b^q and covariance matrix Σ_b^q . From (4'), the mode estimate β^q is updated by

$$\beta^q = \arg \max_{\beta} \{ \exp (E_{-\beta}(\log p(y, \theta))) \} = \arg \max_{\beta} \left\{ \left[\frac{1}{\phi} \right] \int (y' \eta - 1' \zeta(\eta)) q(b) db - \sum_{j=1}^p [\lambda_j] |\beta_j| \right\}. \quad (14)$$

Hereafter, $[\cdot]$ denotes the expectation with respect to the VB posterior. Solving this optimization problem is discussed in detail later on.

For the normal linear mixed regression model, the optimal VB posterior $q(b)$ is a normal distribution and therefore the parameters μ_b^q, Σ_b^q are updated in closed form. In the other cases, from (9), the optimal VB approximation $q(b)$ is

$$q(b) \propto \exp \left(-\frac{1}{2} b' [Q_b] b + \left[\frac{1}{\phi} \right] (y' \eta - 1' \zeta(\eta)) \right) \quad (15)$$

with $\eta = X\beta^q + Zb$. This distribution does not have the form of a standard distribution. We suggest using the Gaussian approximation to approximate this optimal distribution by a normal distribution with mean μ_b^q and covariance matrix Σ_b^q . Let b^* be the maximizer of the function

$$h(b) = -\frac{1}{2} b' [Q_b] b + \left[\frac{1}{\phi} \right] (y' \eta - 1' \zeta(\eta)),$$

which can be easily found by the Newton-Raphson method (see Appendix C). Then, μ_b^q and Σ_b^q are updated as follows

$$\begin{aligned} \mu_b^q &= b^* \\ \Sigma_b^q &= \left(\left[\frac{1}{\phi} \right] Z' \text{diag} \left(\ddot{\zeta}(\eta^*) \right) Z + [Q_b] \right)^{-1}. \end{aligned} \quad (16)$$

with $\eta^* = X\beta^q + Zb^*$.

The optimal VB posterior $q(Q)$ is a Wishart with degrees of freedom and scale matrix

$$\nu^q = \nu_0 + m, \quad S^q = \left(S_0^{-1} + \sum_{i=1}^m (\mu_{b_i}^q \mu_{b_i}^{q'} + \Sigma_{b_i}^q) \right)^{-1}, \quad (17)$$

where $\mu_{b_i}^q$ and $\Sigma_{b_i}^q$ are extracted from μ_b^q and Σ_b^q accordingly. Then, $[Q_b] = \text{diag}([Q], \dots, [Q])$ with $[Q] = \nu^q S^q$.

The optimal VB posterior of λ_j is Gamma with shape and rate

$$\alpha_{\lambda_j}^q = r + 1, \quad \beta_{\lambda_j}^q = |\beta_j^q| + s, \quad (18)$$

and therefore $[\lambda_j] = \alpha_{\lambda_j}^q / \beta_{\lambda_j}^q$. In many cases such as Poisson and logistic regression, ϕ is a known constant, otherwise we can put a suitable prior on ϕ such that the optimal

$$q(\phi) \propto \exp (E_{-\phi}(\log p(y, \theta))) \quad (19)$$

belongs to a recognizable family. In the case of normal linear mixed regression, for example, if using an inverse Gamma prior with shape $\alpha_{\sigma^2}^0$ and scale $\beta_{\sigma^2}^0$ for the dispersion parameter $\phi = \sigma^2$, the optimal VB posterior $q(\sigma^2)$ is an inverse Gamma with shape and scale

$$\alpha_{\sigma^2}^q = n/2 + \alpha_{\sigma^2}^0, \quad \beta_{\sigma^2}^q = \frac{1}{2} \|y - X\beta^q - Z\mu_b^q\|^2 + \frac{1}{2} \text{tr}(Z\Sigma_b^q Z') + \beta_{\sigma^2}^0.$$

In this case, $[1/\sigma^2] = \alpha_{\sigma^2}^q / \beta_{\sigma^2}^q$.

We summarize below the VB algorithm for doing variable selection in GLMMs.

VBGLMM algorithm.

1. Initialize β^q and S^q (and $q(\phi)$ if applicable).
2. Update $\alpha_{\lambda_j}^q$ and $\beta_{\lambda_j}^q$ as in (18).
3. Update μ_b^q and Σ_b^q as in (16)
4. Update S^q as in (17).
5. Update β^q as in (14).
6. Update $q(\phi)$ (if applicable).
7. Repeat Steps 2-6 until convergence.

We may initialize β^q to some initial estimate such as the MLE if available. We suggest to stop the iteration when the difference between two successive updates of the main parameters β^q is smaller than some prespecified value.

Selection of the hyperparameters. For the prior on the λ_j , one can use the improper scale-invariant prior $p(\lambda_j) \propto 1/\lambda_j$, i.e. $r=s=0$. In this paper, we use the empirical Bayes method as in Park and Casella (2008) and Leng et al. (2013) for selecting r . We use a Gamma prior, $\text{Gamma}(\alpha_r^0, \beta_r^0)$, for r and approximating the posterior $p(r|y)$ by $\text{Gamma}(\alpha_r^q, \beta_r^q)$, in which the VB parameters α_r^q, β_r^q are estimated by the fixed-form VB method of Salimans and Knowles (2013). The fixed-form VB algorithm for updating α_r^q, β_r^q is presented in Appendix A. Empirical Bayes update of s is easier, one can put a Gamma prior on s , then the VB optimal posterior of s is also a Gamma. However, we found that, for high-dimensional problems, fixing s to some very small value works better. We set $s = 1e-5$ in our implementation, which implies that we use a very flat prior for the λ_j . We set $S_0 = 10^4 I$ and $\nu_0 = u + 1$ in order to have a flat prior on Q .

3.1 Solving (14)

This section presents a method for solving the optimization problem (14). Let

$$f(\beta) = \left[\frac{1}{\phi}\right] \int (1'\zeta(\eta)) - y'\eta) q(b) db. \quad (20)$$

(14) is equivalent to

$$\arg \min_{\beta} \left\{ F(\beta) = f(\beta) + \sum_{j=1}^p [\lambda_j] |\beta_j| \right\}. \quad (21)$$

It's worth noting that the main different between (21) and (12) is that the integral in $f(\beta)$ can be either computed analytically or approximated easily with an arbitrary accuracy without relying on the Laplace approximation. In (20) we work with the log-scales of the likelihood, which is more convenient than with the original scale as in (12).

Recall that $\eta_{ij} = \beta_0 + x'_{ij}\beta_{1:p} + z'_{ij}b_i$ with $b_i \sim \mathcal{N}(\mu_{b_i}^q, \Sigma_{b_i}^q)$. For normal and Poisson regression $\zeta(\eta_{ij}) = \eta_{ij}^2$ and $\zeta(\eta_{ij}) = e^{\eta_{ij}}$ respectively, the integral in $f(\beta)$ is computed in closed form. After some algebra, it can be shown that

$$f(\beta) = 1' \exp \left(X\beta + Z\mu_b^q + \frac{1}{2} \text{diag}(Z\Sigma_b^q Z') \right) - y'(X\beta + Z\mu_b^q)$$

for Poisson regression. For binomial regression, a closed form approximation to $f(\beta)$ with an arbitrary accuracy is presented in Appendix B.

That is, the function $f(\beta)$ is either computed analytically or easily approximated with an arbitrary accuracy. With a little abuse of notation, we still denote the approximation by $f(\beta)$ in the latter case. An advantage over the method in Groll and Tutz (2012) and Schelldorfer et al. (2013) is that our method does not rely on the Laplace approximation for integrating out the random effects. The Laplace approximation of the likelihood in GLMMs might be in some cases not very accurate (see, e.g. Joe, 2008).

The optimization problem (21) belongs to a popular class of optimization problems in which the target has the form of a sum of a smooth function and a separable convex function. There are many algorithms available for solving such an optimization problem. In this paper, we use the coordinate gradient descent method of Tseng and Yun (2009) (see also Schelldorfer et al., 2013) to solve (21).

Using the notation in Schelldorfer et al. (2013), denote by $\beta^{(s)} = (\beta_0^{(s)}, \dots, \beta_p^{(s)})'$ the value of β at the s th iteration and let $\beta^{(s,s-1;j)} = (\beta_0^{(s)}, \dots, \beta_{j-1}^{(s)}, \beta_j^{(s-1)}, \dots, \beta_p^{(s-1)})'$. Let e_j be the $(j+1)$ st unit vector and $H_j^{(s)}$ be a positive definite matrix, $j=0, \dots, p$. The coordinate gradient descent method is as follows, whose convergence to a stationary point of $F(\beta)$ is proved in Tseng and Yun (2009).

1. Initialize $\beta^{(0)}$. Repeat the following for $s=1, 2, \dots$

2. For $j=0, 1, \dots, p$

(i) Calculate the descent direction

$$d_j^{(s)} = \arg \min_d \left\{ d \nabla f(\beta^{(s,s-1;j)})' e_j + \frac{1}{2} d^2 e_j' H_j^{(s)} e_j + [\lambda_j] |\beta_j^{(s-1)}| + d \right\}. \quad (22)$$

(ii) Choose a step size $\alpha_j^{(s)}$ and set $\beta^{(s,s-1;j+1)} = \beta^{(s,s-1;j)} + \alpha_j^{(s)} d_j^{(s)} e_j$.

For matrix $H_j^{(s)}$ we choose $H_j^{(s)} = \nabla^2 f(\beta^{(s,s-1;j)})$. It is easy to see that $d_j^{(s)}$ in (22) can be solved analytically

$$d_j^{(s)} = \begin{cases} -\frac{\nabla f(\beta^{(s,s-1;j)})' e_j}{e_j' H_j^{(s)} e_j}, & j = 0 \\ \text{median} \left(\frac{[\lambda_j] - \nabla f(\beta^{(s,s-1;j)})' e_j}{e_j' H_j^{(s)} e_j}, -\beta_j^{(s-1)}, \frac{-[\lambda_j] - \nabla f(\beta^{(s,s-1;j)})' e_j}{e_j' H_j^{(s)} e_j} \right), & j > 0. \end{cases}$$

For the step size $\alpha_j^{(s)}$, Tseng and Yun (2009) suggest the Armijo rule as follows: For some $0 < \delta, \varrho < 1$ and $0 \leq \gamma < 1$, choose $\alpha_j^{\text{init}} > 0$ and let $\alpha_j^{(s)}$ be the largest element of $\{\alpha_j^{\text{init}} \delta^l\}_{l=0,1,\dots}$ satisfying

$$F(\beta^{(s,s-1;j)} + \alpha_j^{(s)} d_j^{(s)} e_j) \leq F(\beta^{(s,s-1;j)}) + \alpha_j^{(s)} \varrho \Delta_j,$$

where $\Delta_j = d_j^{(s)} \nabla f(\beta^{(s,s-1;j)})' e_j + \gamma (d_j^{(s)})^2 e_j' H_j^{(s)} e_j$ for $j=0$, and $= d_j^{(s)} \nabla f(\beta^{(s,s-1;j)})' e_j + \gamma (d_j^{(s)})^2 e_j' H_j^{(s)} e_j + [\lambda_j] (|\beta_j^{(s-1)} + d_j^{(s)}| - |\beta_j^{(s-1)}|)$ for $j > 0$. Following Schelldorfer et al. (2013), we choose $\alpha_j^{\text{init}} = 1$, $\delta = 0.5$, $\varrho = 0.1$ and $\gamma = 0$.

4 Examples

4.1 Simulation study

We simulate data sets from a mixed effect Poisson regression model

$$p(y_{ij}|\beta, b_i) = \text{Poisson}(\exp(\eta_{ij})),$$

and a mixed effect logistic regression model

$$p(y_{ij}|\beta, b_i) = \text{Binomial}\left(\frac{\exp(\eta_{ij})}{1 + \exp(\eta_{ij})}\right),$$

with $\eta_{ij} = \beta_0 + x'_{ij}\beta_{1:p} + z'_{ij}b_i$, $i = 1, \dots, n_i$ and $j = 1, \dots, m$. Here, $\beta_0 = 3$ and the first four entries of $\beta_{1:p}$ are $(-2.5, 0, 0, -2)$ and the rest $p-4$ entries are zeros, x_{ij} and z_{ij} are independently generated from the uniform distribution on $(0,1)$, and $b_i \sim \mathcal{N}(0, Q^{-1})$ with $Q = (1/\sigma^2)\mathbb{I}_u$, n_i is set to 5.

We investigate the performance of the proposed VBGLMM approach and compare it to the GLMMLASSO method of Groll and Tutz (2012). We select the best shrinkage parameter λ in the GLMMLASSO method based on BIC from a range of 100 equally-spaced values between 0 and λ_{\max} . Theoretically, λ_{\max} is the smallest value of λ such that $\beta_{1:p} = 0$. Determining λ_{\max} is not straightforward and we set in this simulation example $\lambda_{\max} = 100$ after some experiments.

The performance is measured by the rate of correctly-fitted models (CFR), mean squared errors in β (MSE_β), mean squared errors in σ^2 (MSE_{σ^2}), and CPU time in seconds, over 50 replications.

The simulation results are summarized in Table 1 and Table 2 for various scenario with different values of p , m and σ^2 . VBGLMM outperforms GLMMLASSO in all cases. Especially, VBGLMM works very well in terms of identifying correctly the zero-coefficients.

4.2 Skin cancer data

A clinical trial is conducted to test the effectiveness of beta-carotene in preventing non-melanoma skin cancer (Greenberg et al., 1989). Patients were randomly assigned to a control or treatment group and biopsied once a year to ascertain the number of new skin cancers since the last examination. The response y_{ij} is a count of the number of new skin cancers in year j for the i th subject. The covariates include **age**, **skin** (1 if skin has burns and 0 otherwise), **gender**, **exposure** (a count of the number of previous skin cancers), **year** of follow-up and **treatment** (1 if the subject is in the treatment group and 0 otherwise). There are $m = 1683$ subjects with complete covariate information.

Donohue et al. (2011) argue that **treatment** is not significant and consider 5 different Poisson mixed models with different inclusion of the rest 5 covariates. By using an AIC-type model selection criterion, Donohue et al. (2011) select a random intercept model with four

p	m	σ^2	Method	CFR(%)	MSE_β	MSE_{σ^2}	CPU (seconds)
5	50	0.5	glmmlasso	0	0.123	0.027	123.7
			vbglmm	100	0.091	0.018	3.5
		1	glmmlasso	0	0.140	0.031	278.5
			vbglmm	99	0.101	0.016	5.8
	100	0.5	glmmlasso	0	0.092	0.024	377.7
			vbglmm	100	0.079	0.018	9.1
		1	glmmlasso	0	0.105	0.028	1491.7
			vbglmm	100	0.092	0.022	32.5
50	50	0.5	glmmlasso	0	1.822	0.060	394.7
			vbglmm	85	0.528	0.035	17.3
		1	glmmlasso	0	1.844	0.121	604.6
			vbglmm	81	0.188	0.051	19.5
	100	0.5	glmmlasso	0	0.758	0.038	2226.9
			vbglmm	89	0.481	0.025	44.3
		1	glmmlasso	0	0.738	0.131	941.9
			vbglmm	82	0.291	0.044	17.4

Table 1: Simulation: mixed Poisson regression

p	m	σ^2	Method	CFR(%)	MSE_β	MSE_{σ^2}	CPU (seconds)
5	50	0.5	glmmlasso	0	1.372	0.042	16.3
			vbglmm	98	0.580	0.017	5.6
		1	glmmlasso	0	2.621	0.469	20.3
			vbglmm	89	0.675	0.321	4.7
	100	0.5	glmmlasso	0	1.127	0.055	63.0
			vbglmm	100	0.541	0.015	17.2
		1	glmmlasso	0	1.764	0.521	102.3
			vbglmm	91	0.656	0.189	21.1
50	50	0.5	glmmlasso	0	12.408	0.039	48.9
			vbglmm	72	1.118	0.035	33.2
		1	glmmlasso	0	12.300	0.475	65.6
			vbglmm	72	1.466	0.117	32.1
	100	0.5	glmmlasso	0	5.306	0.067	141.0
			vbglmm	74	0.796	0.056	58.8
		1	glmmlasso	0	5.281	0.554	173.5
			vbglmm	80	1.139	0.157	47.1

Table 2: Simulation: mixed logistic regression

fixed effect covariates `age`, `skin`, `gender`, `exposure` (the fixed effect intercept is always included).

We consider the variable selection problem for this Poisson mixed regression model with a random intercept. We consider all the 6 potential covariates `age`, `skin`, `gender`, `exposure`, `treatment` and `year`. Our method selects the same model as selected by Donohue et al. (2011). The estimate of vector β is $(-24.609, 0.008, 0.350, 1.579, 0.854, 0, 0)$, and the estimate of the random effect standard deviation σ is 102.7.

4.3 Six city data

The six cities dataset in Fitzmaurice and Laird (1993) consists of binary responses y_{ij} which indicate the wheezing status (1 if wheezing, 0 if not wheezing) of the i th child at time-point j , $i=1,\dots,537$ and $j=1,\dots,4$. The covariates are `Age` (the age of the child at time-point j , centered at 9 years) and `Smoke` (the maternal smoking status 0 or 1). We consider the following logistic mixed regression model with two random effects

$$\begin{aligned} p(y_{ij}|\beta, b_i) &= \text{Binomial}(1, p_{ij}), \\ \text{logit}(p_{ij}) &= \beta_0 + \beta_1 \text{Age}_{ij} + \beta_2 \text{Smoke}_{ij} + b_{i1} + b_{i2} \text{Age}_{ij}. \end{aligned}$$

The VBGLMM estimate of β is $(-6.98, 0, 0)$, i.e. `Age` and `Smoke` are not selected. The estimate of the covariance matrix of the random effects b_i is

$$\widehat{\text{Cov}}(b_i) = \begin{pmatrix} 34.863 & -1.103 \\ -1.103 & 0.434 \end{pmatrix}.$$

5 Conclusions and Discussions

We have described in this article a VB algorithm for simultaneous variable selection and parameter estimation in GLMMs. The proposed algorithm is based on the VB method for estimating a posterior mode in conjunction with the Bayesian adaptive Lasso. The posterior mode VB method described in this article can be applied to variable selection in other frameworks such as covariance selection. The proposed VBGLMM method can also be extended to (i) grouped variable selection in GLMMs by using the group lasso penalty (Yuan and Lin, 2006) (ii) ordered variable selection in GLMMs by the composite absolute penalty (Zhao et al., 2009). This research is currently in progress.

Appendix A: Fixed-form VB algorithm for approximating $p(r|y)$

This section presents the fixed-form VB approach of Salimans and Knowles (2013) for approximating $p(r|y)$. Their fixed-form VB algorithm requires an unbiased estimate of a covariance matrix of the form $\text{cov}(T(X), V(X))$ with $T(\cdot)$ and $V(\cdot)$ vector functions of a random variable X with probability density function $f(x)$. Let X_1 and X_2 be two independent draws from f . It is easy to see that

$$\widehat{\text{cov}} = \frac{1}{2}(T(X_1) - T(X_2))(V(X_1) - V(X_2))'$$

is an unbiased estimate of $\text{cov}(T(X), V(X))$.

We use a Gamma prior $\text{Gamma}(\alpha_r^0, \beta_r^0)$ for r and approximate the posterior $p(r|y)$ by $q(r) = \text{Gamma}(\alpha_r^q, \beta_r^q)$. The sufficient statistic for the natural parameter $\eta = (\alpha_r^q, \beta_r^q)'$ is $T(r) = (\log r, -r)'$ and

$$\log p(r, y) = \left(p \log s - \beta_r^0 + \sum_{j=1}^p [\log \lambda_j] \right) r + (\alpha_r^0 - 1) \log r - p \log \Gamma(r)$$

after ignoring the terms independent of r . Let

$$C = C(\alpha_r^q, \beta_r^q) = \begin{pmatrix} \dot{\psi}(\alpha_r^q) & -\frac{1}{\beta_r^q} \\ -\frac{1}{\beta_r^q} & \frac{\alpha_r^q}{\beta_r^{q^2}} \end{pmatrix}.$$

We have the following algorithm for estimating α_r^q and β_r^q .

1. Initialize $\eta = (\alpha_r^q, \beta_r^q)'$. Compute $C = C(\alpha_r^q, \beta_r^q)$ and $g = C\eta$.
2. Initialize $\bar{C} = 0, \bar{g} = 0$.
3. For $i = 1, 2, \dots, N$
 - Set $\eta = C^{-1}g$
 - Generate r_1, r_2 from $q(r)$ and compute

$$\hat{g}_i = \frac{1}{2}(\log p(r_1, y) - \log p(r_2, y))(T(r_1) - T(r_2))$$

and $\hat{C}_i = C(\alpha_r^q, \beta_r^q)$.

- Set $g = (1 - c)g + c\hat{g}_i$, $C = (1 - c)C + c\hat{C}_i$.
 - If $i > N/2$ set $\bar{g} = \bar{g} + \hat{g}_i$, $\bar{C} = \bar{C} + \hat{C}_i$.
4. Set $\eta = \bar{C}^{-1}\bar{g}$.

Appendix B

For binomial mixed regression, $\zeta(\eta_{ij}) = \log(1 + e^{\eta_{ij}})$, where η_{ij} is normally distributed with mean $\beta_0 + x'_{ij}\beta_{1:p} + z'_{ij}\mu_{b_i}^q$ and variance $z'_{ij}\Sigma_{b_i}^q z_{ij}$. The function $f(\beta)$ in (20) becomes

$$f(\beta) = \sum_{i,j} \mathbb{E}_{\eta_{ij}}(\log(1 + e^{\eta_{ij}})) - y'(\beta_0 + X\beta_{1:p} + Z\mu_b^q).$$

Computing $f(\beta)$ reduces to computing the integrals of the form $\mathbb{E}_\xi(\log(1 + e^\xi))$ with $\xi \sim \mathcal{N}(\mu, \sigma^2)$. We write $\mathbb{E}_\xi(\log(1 + e^\xi)) = \mathbb{E}_\zeta h(\zeta)$ with $h(\zeta) = \log(1 + e^{\mu + \sigma\zeta})$ and $\zeta \sim \mathcal{N}(0, 1)$. Using the Taylor expansion of $h(\zeta)$ at zero, $h(\zeta)$ can be approximated by

$$h(\zeta) \approx h(0) + \sum_{k=1}^K \frac{h^{(k)}(0)}{k!} \zeta^k$$

for some $K \geq 1$. Hence,

$$\mathbb{E}_\xi(\log(1 + e^\xi)) \approx h(0) + \sum_{k=1}^K \frac{h^{(k)}(0)}{k!} \mathbb{E}_\xi(\zeta^k).$$

Note that $\mathbb{E}_\zeta(\zeta^k) = 0$ if k is odd and $\mathbb{E}_\zeta(\zeta^k) = (k-1)!!$ if k is even, where $(k-1)!! = 1.3\dots(k-1)$, i.e. the product of every odd number from 1 to $k-1$. We set $K=2$ in the examples reported in this article. The user can set a bigger K in the R package `vbglm`.

Appendix C: Gaussian approximation

Suppose that $q(x) = e^{f(x)}$ and we wish to approximate $q(x)$ by a Gaussian density. Let x^* be the maximizer of $f(x)$. By Taylor's expansion

$$f(x) \approx f(x^*) + \frac{1}{2}(x - x^*)' \frac{\partial^2 f(x^*)}{\partial x \partial x'} (x - x^*).$$

Then

$$\begin{aligned} q(x) = e^{f(x)} &\approx \exp\left(f(x^*) + \frac{1}{2}(x - x^*)' \frac{\partial^2 f(x^*)}{\partial x \partial x'} (x - x^*)\right) \\ &\propto \exp\left(\frac{1}{2}(x - x^*)' \frac{\partial^2 f(x^*)}{\partial x \partial x'} (x - x^*)\right). \end{aligned}$$

So the best Gaussian approximation to $q(x)$ has mean x^* and covariance matrix $-(\frac{\partial^2 f(x^*)}{\partial x \partial x'})^{-1}$.

Recall that we wish to maximize

$$h(b) = -\frac{1}{2}b'[Q_b]b + \left[\frac{1}{\phi}\right](y'\eta - 1'\zeta(\eta)),$$

with $\eta = X\beta^q + Zb$. The first and second derivatives are

$$\begin{aligned} u(b) = \frac{\partial h(b)}{\partial b} &= \left[\frac{1}{\phi}\right]Z'(y - \zeta(\eta)) - [Q_b]b \\ H(b) = \frac{\partial^2 h(b)}{\partial b \partial b'} &= -\left[\frac{1}{\phi}\right]Z' \text{diag}\left(\ddot{\zeta}(\eta)\right)Z - [Q_b]. \end{aligned}$$

The Newton-Raphson method for maximizing $h(b)$:

1. Initialize b^{old} .
2. Update until some stopping rule is satisfied

$$b^{\text{new}} = b^{\text{old}} - H(b^{\text{old}})^{-1}u(b^{\text{old}}).$$

References

- Dempster, A., Laird, N., and Rubin, D. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 39(1):1–38.
- Donohue, M. C., Overholser, R., Xu, R., and Vaida, F. (2011). Conditional Akaike information under generalized linear and proportional hazards mixed models. *Biometrika*, 98:685–700.
- Fitzmaurice, G. and Laird, N. (1993). A likelihood-based method for analysing longitudinal binary responses. *Biometrika*, 80:141–151.
- Greenberg, E. R., Baron, J. A., Stevens, M. M., Stukel, T. A., Mandel, J. S., Spencer, S. K., Elias, P. M., Lowe, N., Nierenberg, D. N., G., B., and Vance, J. C. (1989). The skin cancer prevention study: design of a clinical trial of beta-carotene among persons at high risk for nonmelanoma skin cancer. *Controlled Clinical Trials*, 10:153–166.
- Griffin, J. E. and Brown, P. J. (2011). Bayesian adaptive Lasso with non-convex penalization. *Australian and New Zealand Journal of Statistics*, 53:423–442.
- Groll, A. and Tutz, G. (2012). Variable selection for generalized linear mixed models by l1-penalized estimation. *Statistics and Computing*, pages 1–18.
- Joe, H. (2008). Accuracy of Laplace approximation for discrete response mixed models. *Computational Statistics & Data Analysis*, 52(12):5066 – 5074.
- Leng, C., Tran, M.-N., and Nott, D. J. (2013). Bayesian adaptive lasso. *The Annals of the Institute of Statistical Mathematics*. To appear.
- Park, T. and Casella, G. (2008). The bayesian lasso. *Journal of the American Statistical Association*, 103:681–686.
- Salimans, T. and Knowles, D. A. (2013). Fixed-form variational posterior approximation through stochastic linear regression. Technical report, Erasmus University Rotterdam. Available at <http://arxiv.org/abs/1206.6679>.
- Schelldorfer, J., Meier, L., and Bhlmann, P. (2013). GLMMLasso: An algorithm for high-dimensional generalized linear mixed models using l1-penalization. *Journal of Computational and Graphical Statistics*, 0(ja):null.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society B*, 58(1):267–288.
- Tseng, P. and Yun, S. (2009). A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1-2):387–423.
- Yuan, M. and Lin, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society, Series B*, 68:49–67.
- Zhao, P., Rocha, G., and Yu, B. (2009). The composite absolute penalties family for grouped and hierarchical variable selection. *The Annals of Statistics*, 37:34683497.

Zou, H. (2006). The adaptive Lasso and its oracle properties. *Journal of the American Statistical Association*, 101(476):1418–1429.